# Using Graphs
# to Analyze Big Linked Data

**Hassan Chafi,**
**Director, Research and Advanced Development**
**Oracle Labs**

**ORACLE**

# Safe Harbor Statement

The following is intended to outline our research activities and general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

# Big Data and Data Analysis

- The Big Data era is here
  - Volume
  - Velocity
  - Variety

- However, just storing and managing this data is not sufficient
  - Typically Big Data is low value per byte

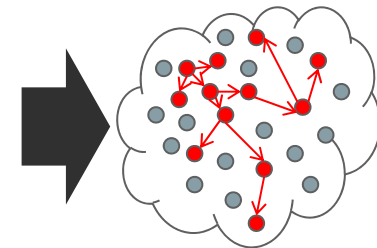- ➔ We want to get useful information out of the huge data sets

Data Analytics:
- Classic OLAP/BI
- Statistical analysis
- Machine learning
- Graph analysis

ORACLE

# Graph Analysis

- Represent your data as a graph
  - Data entities become vertices
  - Relationships become edges



- Analyzing the graph can yield very useful information
  - As it inspects fine-grained relationships between data entities and can factor in the shape of the resulting graph

- Challenge:
  - Traditional tools are neither convenient nor efficient for graph analysis
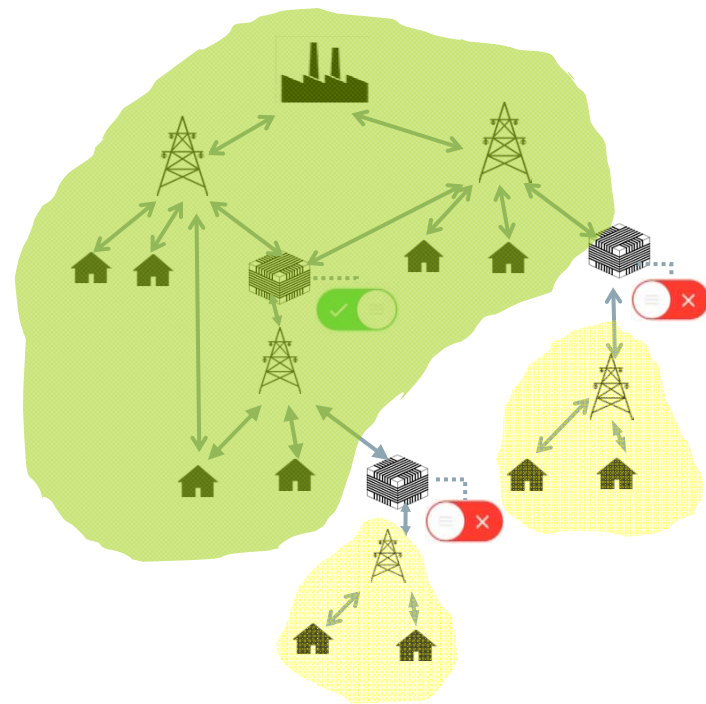
➔ Let's start with a very simple real world example

# Example #1: Power Distribution Network

- Real world use case from an infrastructure company

- [The data set]
  - The dataset represents a power distribution network
  - Generators, transformers, lines, regulators, switches, …

- [One simple question] (out of many)
  - What/how many elements still remain reachable from a source, when given switches are turned off

ORACLE

# Example #1: Relational Data Model

- Their data was originally stored in relational tables*

| ID | Name | Type | Voltage... |
|---|---|---|---|
| 2018281 | XFM_Sub | Generator | ... |
| 27080172 | SW_35 | Switch | ... |
| ... | | | |

Table for node information

| Connection | Node A | Node B | Node C | Node D | Node E... |
|---|---|---|---|---|---|
| 6693 | 2018281 | 289301 | 4985701 | - | - |
| 7207 | 2019182 | 495812 | 9191913 | 4985701 | - |
| ... | ... | ... | ... | | |

Table for connection information (using null values for variable connection size)

*Schema simplified for explanation's sake

**Issue #1**
The question cannot be answered with simple SQL
➔ Cannot avoid some programming (e.g. PL/SQL)

**Issue #2**
Performance issues (next slide)

ORACLE®

# Example #1: Performance Issue

- What needs to be done:

Overall, many joins during each of many iterations!

| ID | Name | Type | Voltage... |
|---|---|---|---|
| 2018281 | XFM_Sub | Generator | ... |
| 27080172 | SW_35 | Switch | ... |
| ... | | | |

Need to exclude switches that are off

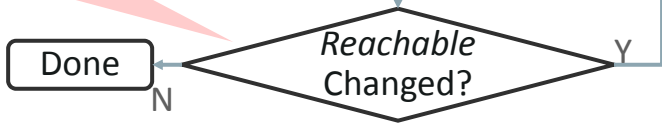| Connection | Node A | Node B | Node C | Node D | Node E... |
|---|---|---|---|---|---|
| 6693 | 2018281 | 289301 | 4985701 | - | - |
| 7207 | 2019182 | 495812 | 9191913 | 4985701 | - |
| ... | ... | ... | ... | ... | |
| | | | 289301 | | |

Repeat until not changed

Find start node from node table.
*Reachable* = {start_node}
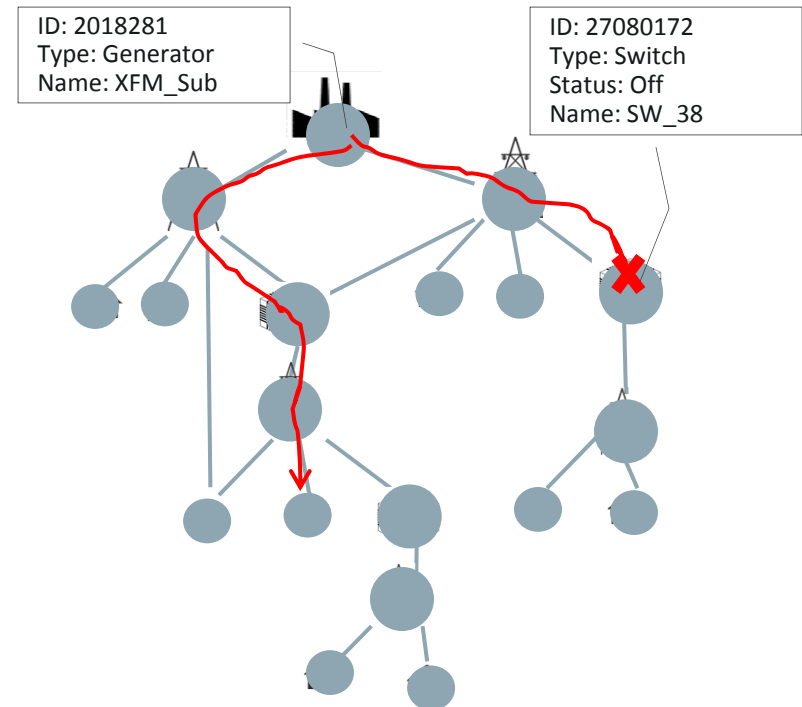
Find *Next* reachable nodes by joining connection table

*Reachable* = *Reachable* U *Next*

*Reachable* Changed?  — Y

Done — N

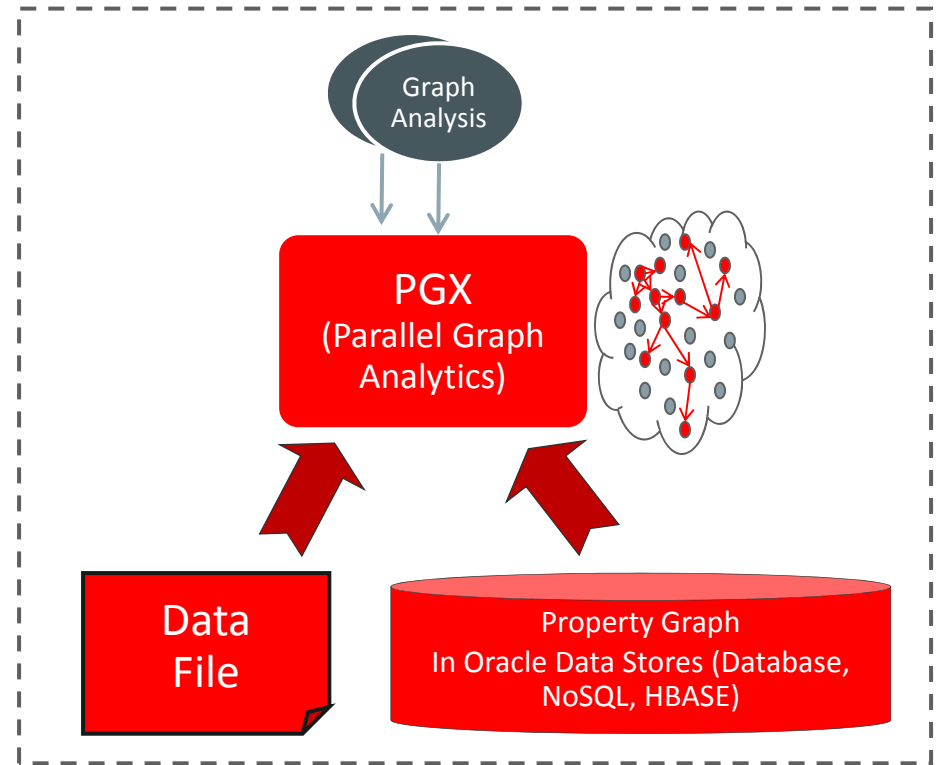# Example #1: Use A Graph Model, Instead

- Represent the data as a graph
  - Vertices and edges have extra information, or *properties*
  - Fits very naturally

- Answer the question in natural ways
  - Starting from the given vertex,
  - traverse the graph,
  - without going through 'off' switches

> But, what tools provide such data model and operations?

ID: 2018281
Type: Generator
Name: XFM_Sub

ID: 27080172
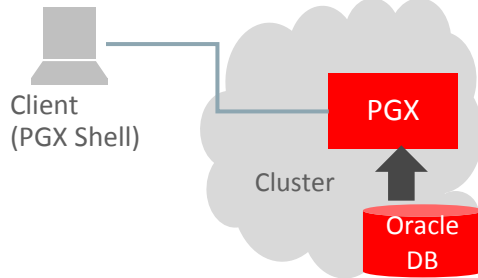Type: Switch
Status: Off
Name: SW_38

# Introducing PGX

- PGX is a graph processing system from Oracle Labs

- PGX reads graph data from files or Oracle DBs

- PGX builds and maintains its own graph representation

- … and provides graph analysis features to the users



Graph
Analysis

PGX
(Parallel Graph
Analytics)

Data
File

Property Graph
In Oracle Data Stores (Database,
NoSQL, HBASE)

**ORACLE**

# Example #1: PGX Example



Client
(PGX Shell)

PGX

Cluster

Oracle
DB

- PGX supports server-client mode execution
- Multiple concurrent clients
- *Groovy*-(and iPython) based interactive shell

**Load graph and get root vertex**

```
pgx> G = session.readGraphWithProperties(" … ")
pgx> root = G.getVerteX(21474836490)
```

**Create a filter condition.**
**Do BFS(Breadth-First Search) traversal from root**

```
pgx> navigator = new VertexFilter(
    "vertex.type!='switch' or vertex.state==true")
pgx> analyst.filteredBfs( G, root, null, navigator )
```

**Find reached nodes, and print their distances from the root**

```
pgx> G.queryPgql (
    "SELECT vertex, vertex.distance    \
    WHERE (vertex WITH distance >=0) \
    ORDER BY vertex.distance DESC").print(30)
```

ORACLE®

# Example #2: Influencer Identification

- General Idea
  - Capture relationships between data entities as a graph
  - Inspect the topology of the graph and identify *important* entities

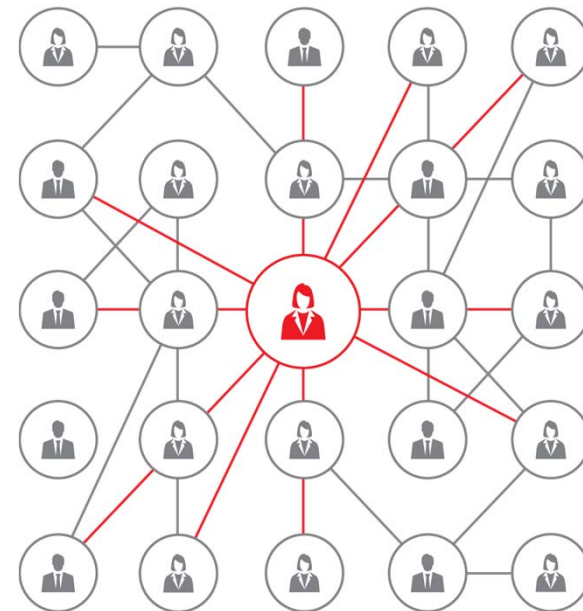- More Theoretic Approach
  - Centrality in Graph theory

| Article | Talk | | Read | Edi |
| --- | --- | --- | --- | --- |

## Centrality

From Wikipedia, the free encyclopedia
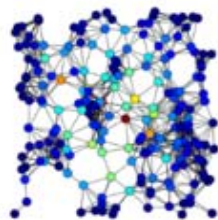
For the statistical concept, see *Central tendency*.

In graph theory and network analysis, indicators of **centrality** identify the most important vertices within a graph. Applications include identifying the most influential person(s) in a
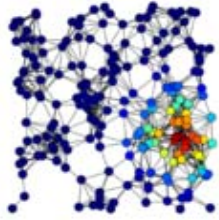
**ORACLE**

# Example #2: Centrality Algorithms?

- Graph theory defines many different measures
  - Betweenness Centrality
  - Closeness Centrality
  - Eigenvector Centrality
  - Pagerank
  - HITS
  - …

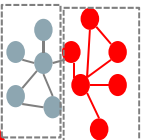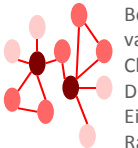Each algorithm suggests different definition of *importance*

Betweenness centrality    Eigenvector centrality

(images from Wikipedia)

- PGX comes with built-in algorithms – various centralities included

**Detecting Components and Communities**

Tarjan's, Kosaraju's, Weakly Connected Components, Label Propagation (w/ variants), Soman and Narang's Spacification

**Ranking and Walking**

Pagerank, Personalized Pagerank, Betwenness Centrality (w/ variants), Closeness Centrality, Degree Centrality, Eigenvector Centrality, HITS, Random walking and sampling (w/ variants)

**Evaluating Community Structures**

Conductance, Modularity Clustering Coefficient (Triangle Counting) Adamic-Adar

**Path-Finding**

Hop-Distance (BFS) Dijkstra's, Bi-directional Dijkstra's Bellman-Ford's

**Link Prediction** SALSA (Twitter's Who-to-follow)

**Other Classics** Vertex Cover Minimum Spanning-Tree (Prim's)

**ORACLE**

# Example #2: Let's try with a *real* data set

- A network of characters in Marvel comics

- http://exposedata.com/marvel/

  1. Download .csv file; it is a simple

     list of *linked* characters

  ```
  "NOVA/RICHARD RIDER","HULK/DR. ROBERT BRUC"
  "NOVA/RICHARD RIDER","NIGHT THRASHER/DUANE"
  "NOVA/RICHARD RIDER","SPIDER-MAN/PETER PAR"
  "SPIDER-MAN/PETER PAR","FIRESTAR/ANGELICA JO"
  "SPIDER-MAN/PETER PAR","THUNDERBALL/DR. ELIO"
  ```

2. Write up a simple JSON descriptor

```
{
 "uri": "hero-network.csv",
 "format": "edge_list",
 "node_id_type": "string",
 "separator": ","
}
```

We want the graph *undirected*

3. Start PGX and load the graph

```
pgx> G = session.readGraphWithProperties("hero-network.csv.json").undirect();
```

# Example #2: Running Centrality Algorithms

```
pgx> analyst.pagerank(G, 0.0001, 0.85, 100)
pgx> analyst.vertexBetweenessCentrality(G)


pgx> G.queryPGQL("\
    SELECT n, n.pagerank, n.betwenness \
    WHERE (n) ORDER BY DESC (n.pagerank)").print(10);



pgx> G.queryPGQL("\
    SELECT n, n.pagerank, n.betweenness \
    WHERE (n) ORDER BY DESC (n.betweenness)").print(10);
```

Compute pagerank and betweenness centrality

Query the graph and get top 10 pagerank and betwenness centrality

Want to see the results? (next slide)

**ORACLE®**

# Two Kinds of Graph Workloads

## Computational Graph Analytics

Connected Components

Modularity    Conductance

Shortest Path

Pagerank    Spanning Tree

Clustering Coefficient

Centrality

Coloring

**Compute** certain values on nodes and edges

While (repeatedly) **traversing** or **iterating** on the graph

In certain **procedural** ways

## Graph Pattern Matching



Given a **description** of a pattern

Find every sub-graph that **matches** it

PGX supports <u>**both**</u> kinds, as well as **combinations of the two**

Images from IMDB.com

# Graph Pattern Matching

- PGQL
  - A query language for graph pattern match
  - SQL-like syntax but with graph pattern description and property access

Example> Retrieve all the neighbors of Captain America, list them ordered by their Pagerank values

```
SELECT n2
WHERE
    (n1 WITH id() = 'CAPTAIN AMERICA') -> n2
ORDER BY n2.pagerank
```

- SQL Like syntax format:
  - SELECT, WHERE, ORDER BY
- Graph pattern description:
  - (n1 … ) -> n2
- Property access :
  - (n1 WITH id() = '…' )
  - n2.pagerank

n1 → n2

Id: 'CAPTAIN AMERICA'

Id: ???

# Example #2: More Information with Patterns

- PGQL enables you to ask more interesting queries

```
pgx> G = G.simplify( … )

pgx> G.queryPgql ("SELECT x \
      WHERE (a WITH id()='SHANG-CHI')->x,\
            (b @ 'WHITE TIGER/HECTOR A')->x,
            (c @ 'IRON FIRST/DANIEL RAN')->x\
").print(20);

pgx> G.queryPgql ("SELECT b, x  \
      WHERE (a@'DR. OCTOPUS/OTTO OCT')->x<-b,\
            x.pagerank < 0.0001 ,\
            b.pagerank > 0.005 \
").print(10);
```

Remove multiple edges to avoid repeated answers

Who are the common neighbors of three martial-art heroes:

Short-cut syntax for Id() matching

Find Dr. Ock's *minor* neighbors which has link to a *major* character (assumed by pagerank value)

Computational analysis combined with pattern matching

IRON MAN
SPIDER-MAN
MOON KNIGHT
FU, MANCHU

**...** (21 results)

SPIDER-MAN
CHANCE II
HUMAN TORCH
GAYLE WATSON

**...** (100 results)

# What if I need other algorithms?

- Graph algorithms are evolving
  - e.g., Scientists keep proposing different definitions of centrality

- What if I want a custom algorithm than built-in one?

- In other words, how programmable is PGX?
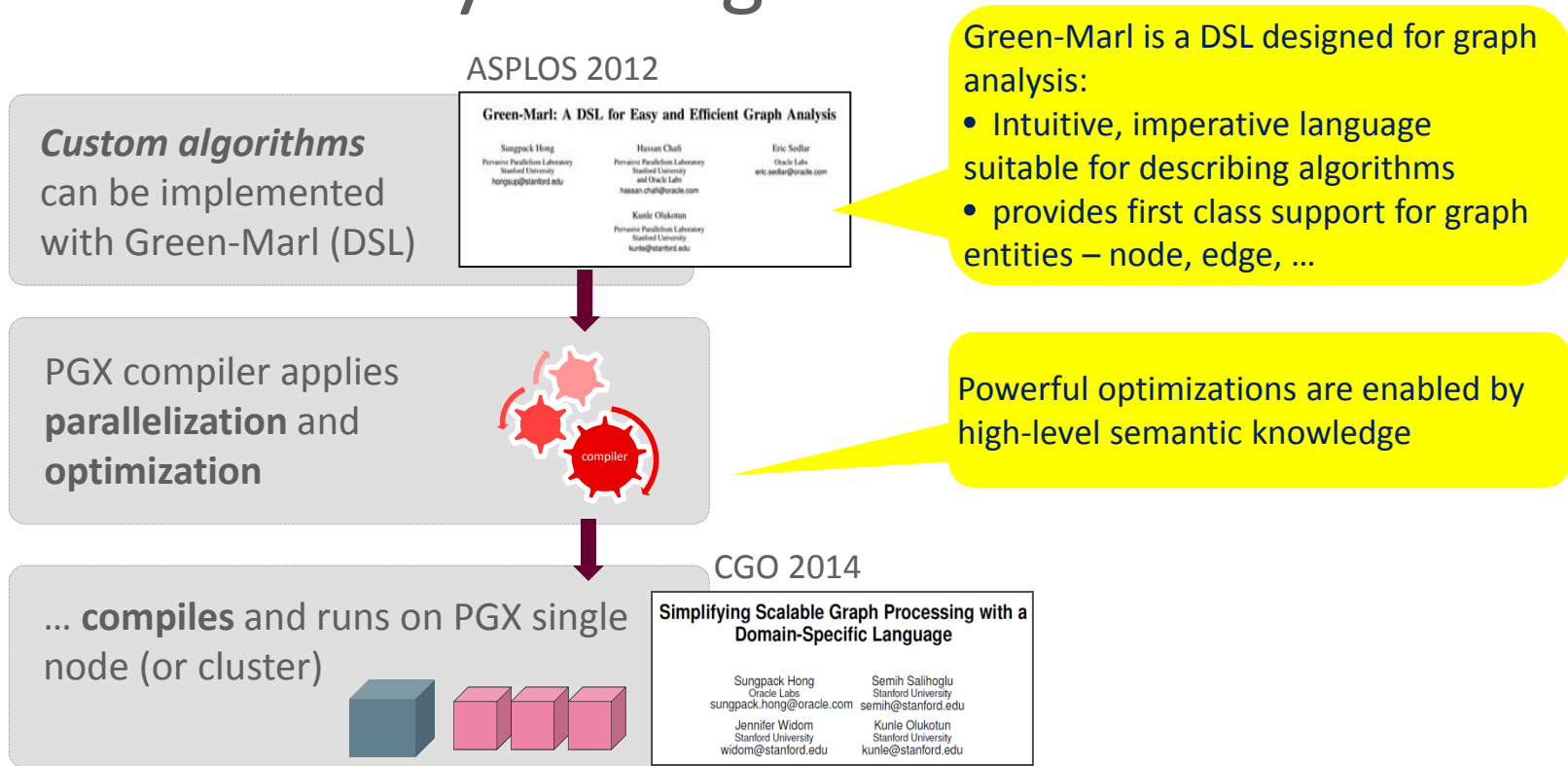
- In two ways:

**(1) Via Java API**
- `getVertex(), getProperty(), …`
- allows fine-grained control
- the user implements his/her algorithm

(2) **Via a Domain-Specific Language (DSL)**
- intuitive programming
- automated parallelization and optimization

(more in next slide)

**ORACLE**

# PGX -- Flexibility through DSL

ASPLOS 2012

**Green-Marl: A DSL for Easy and Efficient Graph Analysis**

Sungpack Hong
Pervasive Parallelism Laboratory
Stanford University
hongsup@stanford.edu

Hassan Chafi
Pervasive Parallelism Laboratory
Stanford University
and Oracle Labs
hassan.chafi@oracle.com

Eric Sedlar
Oracle Labs
eric.sedlar@oracle.com

Kunle Olukotun
Pervasive Parallelism Laboratory
Stanford University
kunle@stanford.edu

***Custom algorithms*** can be implemented with Green-Marl (DSL)

Green-Marl is a DSL designed for graph analysis:
• Intuitive, imperative language suitable for describing algorithms
• provides first class support for graph entities – node, edge, …

PGX compiler applies **parallelization** and **optimization**

compiler

Powerful optimizations are enabled by high-level semantic knowledge

CGO 2014

**Simplifying Scalable Graph Processing with a Domain-Specific Language**

Sungpack Hong
Oracle Labs
sungpack.hong@oracle.com

Semih Salihoglu
Stanford University
semih@stanford.edu

Jennifer Widom
Stanford University
widom@stanford.edu

Kunle Olukotun
Stanford University
kunle@stanford.edu

… **compiles** and runs on PGX single node (or cluster)

ORACLE

# DSL Example: Betweenness Centrality

```
procedure my_bc(G: graph;          // input graph
       BC: nodePropery<double>)   // output node property
{
   G.BC = 0;  // Initialize output

   foreach (s:G.nodes) { // from each node s in G

       // create temporary node properties and initialize them
       nodeProperty<double> sigma;
       nodeProperty<double> delta;
       G.sigma = 0;
       s.sigma = 1;

       // Do BFS traversal from s
       inBFS(v: G.nodes from s) (v != s) {
           // compute sigma for v, by summing over BFS parents
           v.sigma = sum(w:v.upNbrs) { w.sigma };
       }
       inReverse(v!=s) { // Do reverse-BFS order iteration to s
           v.delta =        // compute delta for v, by summing over BFS children
               sum (w:v.downNbrs){(1+ w.delta)/w.sigma} * v.sigma;

           v.BC += v.delta ; // accumulate delta into BC
       }
   }
}
```

**Green-Marl Program**

Writing graph algorithm becomes straight-forward. The language provides graph-specific intrinsics such as graph, property, BFS, neighbors ..

Our intelligent compiler sees parallelizable operations in the algorithm, and (selectively) applies them automatically.

ORACLE

# DSL Example: Betweenness Centrality

```
procedure my_bc(G: graph;          // input graph
        BC: nodePropery<double>)   // output node property
{
  G.BC = 0; // Initialize output

  foreach (s:G.nodes) { // from each nodes s in G

    // create temporary node properties and initialize them
    nodeProperty<double> sigma;
    nodeProperty<double> delta;
    G.sigma = 0;
    s.sigma = 1;

    // Do BFS traversal from s
    inBFS(v: G.nodes from s) (v != s) {
      // compute sigma for v, by summing over BFS parents
      v.sigma = sum(w: v.upNbrs) { w.sigma };
```
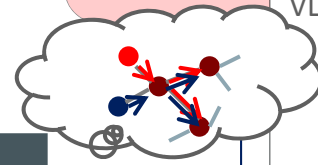
The compiler can apply even more powerful optimizations.

For instance, the compiler sees that BFS is repeated inside a parallel loop

Then the compiler merges multiple BFS traversals into one

VLDB 2015
**The More the Merrier:**
**Efficient Multi-Source Graph Traversal**

Manuel Then[*]          Moritz Kaufmann[*]          Fernando Chirigati[†]          Tuan-Anh Hoang-Vu[†]
then@in.tum.de          kaufmanm@in.tum.de          fchirigati@nyu.edu

Kien Pham[†]          Alfons Kemper[*]          Thomas Neumann[*]          Huy T. Vo[†]
kien.pham@nyu.edu          kemper@in.tum.de          neumann@in.tum.de          huy.vo@nyu.edu

[*] Technische Universität München          [†] New York University

| | PGX | iGraph (C++ library) |
|---|---|---|
| Betwenness Centrality | 80 mins | Not finished after 48 hours |
| Closeness Centrality | 9 mins | Not finished after 48 hours |

Performance is greatly improved !
(measured with 200 million edged graph)

}

Green-Man Program

# How does it work with in PGX?

```
pgx> my_bc = session.compile("my_bc.gm")
```

Read a graph and run the algorithm

```
pgx> G = session.readGraphWithProperties(" … ")
pgx> P = G.createVertexProperty(PropertyType.DOUBLE, "my_prop");
pgx> my_bc.run(G, P);
```

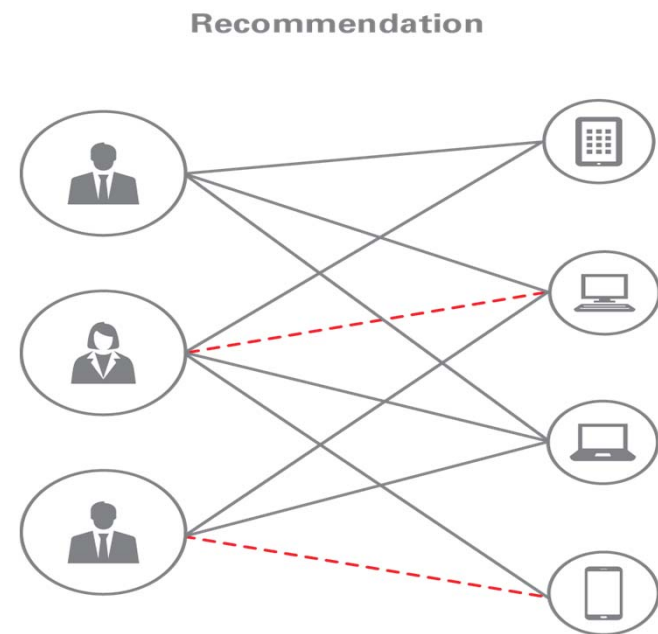Creating a new vertex property to hold the output

Compiled methods can be naturally invoked from Groovy (or Java)

Compare to Built-in Implementation

```
pgx> anayst.vertexBetweennessCentrality(G)
pgx> G.queryPgql (
  "SELECT n, n.my_prop, n.betweenness   \
   WHERE (n) \
   ORDER BY DESC(n.betweenness").print(10);
```

| n | n.my_prop | n.betweenness |
|---|---|---|
| SPIDER-MAN/PETER PAR | 5033505.395979839 | 5033505.395979841 |
| CAPTAIN AMERICA | 3720705.946340431 | 3720705.9463404035 |
| IRON MAN/TONY STARK | 2061886.2266715805 | 2061886.2266715826 |
| WOLVERINE/LOGAN | 2055461.3906177846 | 2055461.3906177809 |
| DR. STRANGE/STEPHEN | 1628817.0459632506 | 1628817.0459632422 |
| HAVOK/ALEX SUMMERS | 1612813.882439691 | 1612813.8824396855 |

# Example #3: Recommendation

- Data set and Question
  - Users and Items
  - Items purchases by each user and their ratings

- [Q] Given a user, can we suggest an item that he/she would like to buy next?

- General Idea: Collaborative Filtering
  - Find users that are *similar* to the given user
  - Find items that are favored by those users
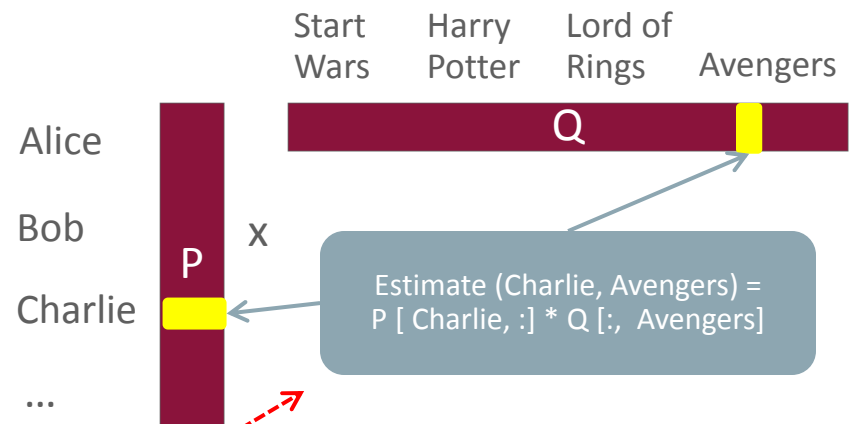  - Recommend such items

**Recommendation**

**ORACLE**

# Example #3: Mathematical Backgrounds

- Matrix Factorization
  - Consider the ratings data as a (sparse) matrix



|  | Start Wars | Harry Potter | Lord of Rings | Avengers | ... |
|---|---|---|---|---|---|
| Alice | ★★★★ | - | ★★ | ★★★★★ | |
| Bob | ★ | ★★★ | ★★★★★ | - | |
| Charlie | ★★★★ | ★★ | ★ | ? | |
| ... | | | | | |

User- Item Matrix : [N x M]
N and M are big Numbers

What would be Charlie's rating for this item?

|  | Start Wars | Harry Potter | Lord of Rings | Avengers |
|---|---|---|---|---|
| Alice | | Q | | |
| Bob | P | x | | |
| Charlie | | | | |
| ... | | | | |

Estimate (Charlie, Avengers) = P [ Charlie, :] * Q [:, Avengers]

Approximate the matrix as multiplication of two matrices: [N x K] * [K x M]

ORACLE®

# Example #3: Graph Intuition

- The exactly same technique can be applied with graph representation
    - The sparse ratings matrix ➔ Bipartite Graph

**We can think that each person or item is characterized as a K-length *feature vector***

**Ratings are the inner product of the two features, or *similarities* between the customer and the item**

[0.758 0.331 0.124 …]

[-0.305 0.888 0.931 ….]

[0.75…

[0.39…

**Customer 's features are defined from the features of customers who bought them.**

**Customer 's features are defined from the features of their purchased items**

0.172 0.519 …]

**Need to solve these recursive equations**

$$\vec{P}_j \;=\; f\left(\sum_{i \in nbr \;(\,j\,)} w_{ij} \cdot \vec{Q}_i\right)$$

$$\vec{Q}_i \;=\; g\left(\sum_{j \in nbr \;(\,i\,)} w_{ij} \cdot \vec{P}_j\right)$$

# Example #3: PGX Application Flow

**1. Data Acquisition**
- Ratings dataset are loaded as a bipartite graph



customer    items

Ratings Dataset

**2. Training**
- Compute feature vectors
- PGX provides matrix factorization implementation

```
while (i < max_step) {
    double rmse = 0.0;

    foreach (u: G.nodes) {
        if (u.is_left) {
            vect<double>[K] Z = 0.0;

            for(e: u.outEdges) {
                node v = e.toNode();

                double w1 = e.weight;
                double w2 = u.PQ * v.PQ;

                if (w2 > VALUE_MAX) {
                    w2 = VALUE_MAX;
                } else if (w2 < VALUE_MIN) {
                    w2 = VALUE_MIN;
                }

                Z += ((w1 - w2) * v.PQ - lambda
                rmse += (w1 - w2) * (w1 - w2);
            }
```

**3. Recommend**
- Given a user and set of items, compute inner products of feature vectors.

➔ Recommend Top-K items

Vermont Is For Lovers - 5.920148497542082
Night Flier - 5.691110431206573
Boys Life - 5.682416396766282
Frisk - 5.654641389008914
Visitors, The (Visiteurs, Les) - 5.579409110367704
Jerky Boys, The - 5.546784101675797
Quartier Mozart - 5.531018328773964
Squeeze - 5.516017286422173
I Can't Sleep (J'ai pas sommeil) - 5.5096872333124285
Crossfire - 5.49346348095436910

# Other Graph Analysis Examples Include …


Shortest path


Anomaly Detection

Community Detection and Evaluation, Link Prediction, Intruder Trace, Network Flow Analysis ……

ORACLE

# PGX: Comparison of Features

| | Graph Database (**Neo4J**, Titan, … ) | Graph Analytics Framework (**GraphX**, GraphLab, …) | PGX |
|---|---|---|---|
| **Parallel or Distributed Graph Analysis** | ✗ | ✓ | ✓ |
| **Pattern Matching (Graph Query)** | ✓ | ✗ | ✓ |
| **Transactional Data Management** | ✓ | ✗ | ✓ |

- Orders of magnitude better performance
- Intuitive DSL and **auto-parallelizing** compiler

- Orders of magnitude better performance
- Expressive graph query language

- Integration with proven Oracle Database (RDBMS, NoSQL, …)

# PGX: Performance Compared with Neo4J

Path queries of Linux kernel source code

```
X86 Server
Xeon E5-2660 2.2Ghz
 2 socket
 x 8 cores
 x 2HT
256GB DRAM

Neo4J: 2.2.1
Data:
 - Linux kernel code as a
graph
 - Program analysis queries
```

## Linux Kernel analysis on X86

■ PGX  ■ Neo4j

**Huge performance advantage over Neo4J graph DB (2~4 orders of magnitude)**

11987x

210x   276x   285x   276x   225x   909x   525x

Time (ms): 100000, 10000, 1000, 100, 10, 1, 0.1, 0.01

Q1   Q2   Q3   Q4   Q5   Q6   Q7   Q8

**Basic graph pattern**

**Path queries**

**Single shortest path**

**Bulk shortest path**

ORACLE®

# Performance Comparison with GraphX, GraphLab

Benchmarking Result
with popular graph
algorithms

**Y-axis: Relative Performance (Higher is Better, Log-scale)**

**PGX outperforms GraphLab and GraphX by orders of magnitude**

```
Cluster
• Xeon E5-2660,2.2Ghz
  2 socket,8 cores
  2HT, 256GB DRAM)
• 2 ~ 32 Machines
• InfiniBand : 56Gbps


DataSet
TWT: 1.4B edges
WEB: 3.0B edges

GraphX:1.1.0
GraphLab: 2.1
```

**X-axis: Number of machines (2 ~ 32)**

**WCC algorithm applied on WEB graph**

**Similar performance gaps across many different algorithms**



Legend:
- PGX (single node)
- PGX (cluster)
- GraphLab
- GraphX

WCC(TWT)

WCC(WEB)

PR(TWT)  PR(WEB)  PR'(TWT)  PR'(WEB)  WCC(TWT)

WCC(WEB)  SSSP(TWT)  SSSP(WEB)  HopDist(TWT)  HopDist(WEB)

EV(TWT)  EV(WEB)  KCORE(LJ)  KCORE(WIK)  Legend

Legend:
- Single Machine Pull
- Single Machine Push
- PGX.D Pull
- PGX.D Push
- GraphLab
- GraphX

ORACLE

# Try PGX Today

## OTN Technology Preview



http://bit.ly/pgxdld

## Oracle Big Data Spatial & Graph

# Architecture of Property Graph Support

**Graph Analytics**

**Parallel In-Memory Graph Analytics (PGX)**

↕ **Java APIs**

**Graph Data Access Layer (DAL)**

**Blueprints & Lucene/SolrCloud**

↕ **Java APIs/JDBC/SQL/PLSQL**

**Scalable and Persistent Storage Management**

**Oracle RDBMS**    **Apache HBase**    **Oracle NoSQL Database**

**Java, Groovy, Python, …**
**REST/Web Service**

**Property Graph formats**

**GraphML**
**GML**
**Graph-SON**
**Flat Files**

# Property Graph Database for Oracle Database 12c Release 2

**NEW IN 12.2**

## Oracle Spatial and Graph

- Neo4J/TITAN-like capabilities in Oracle DB

- Massively-Scalable Graph Database
  - Scales securely to **trillions** edges

- In-Memory and In-Database Parallel Graph Analytics
  - **30+ built-in memory-based parallel** graph analysis algorithms
  - **Optimized in-database** scalable graph algorithms

- Simple interfaces
  - **Parallel SQL queries** for graph analysis and filtering of subgraphs
  - Java: Tinkerpop: Blueprints, Gremlin, Rexster
  - Oracle Text, Apache Lucene and SolrCloud
  - Groovy, Python

**ORACLE**

# Conclusion

- Graph Analysis
  - Inspects the fine-grained relationships between data entities
  - Discover implicit information from the data set
  - Many applications
    - reachability and shortest path, influencer identification, recommendation, community detection, …

- PGX is an easy and efficient tool for graph analysis

- Supported component in Big Data Spatial and Graph and Upcoming 12cR2 release of Spatial and Graph Database Option

# Hardware and Software

## Engineered to Work Together

ORACLE®